

Поиск музыки в социальной сети по запросу с орфографическими ошибками

Ерзикова Юлия
Галкина Алёна
Филатова Анастасия

Райский Андрей
Перваков Григорий
Зархидзе Никита

Группа 341

Постановка задачи

Поиск музыкальных произведений по запросам с орфографическими ошибками и опечатками



Сбор данных



- База данных наиболее запрашиваемых англоязычных песен с сайта lyrics.net (на данный момент lyrics.com)

<u>wc_lyricsnet_albums</u>	
id	PK
lyricsnet_id	
artist_id	
name	
year	

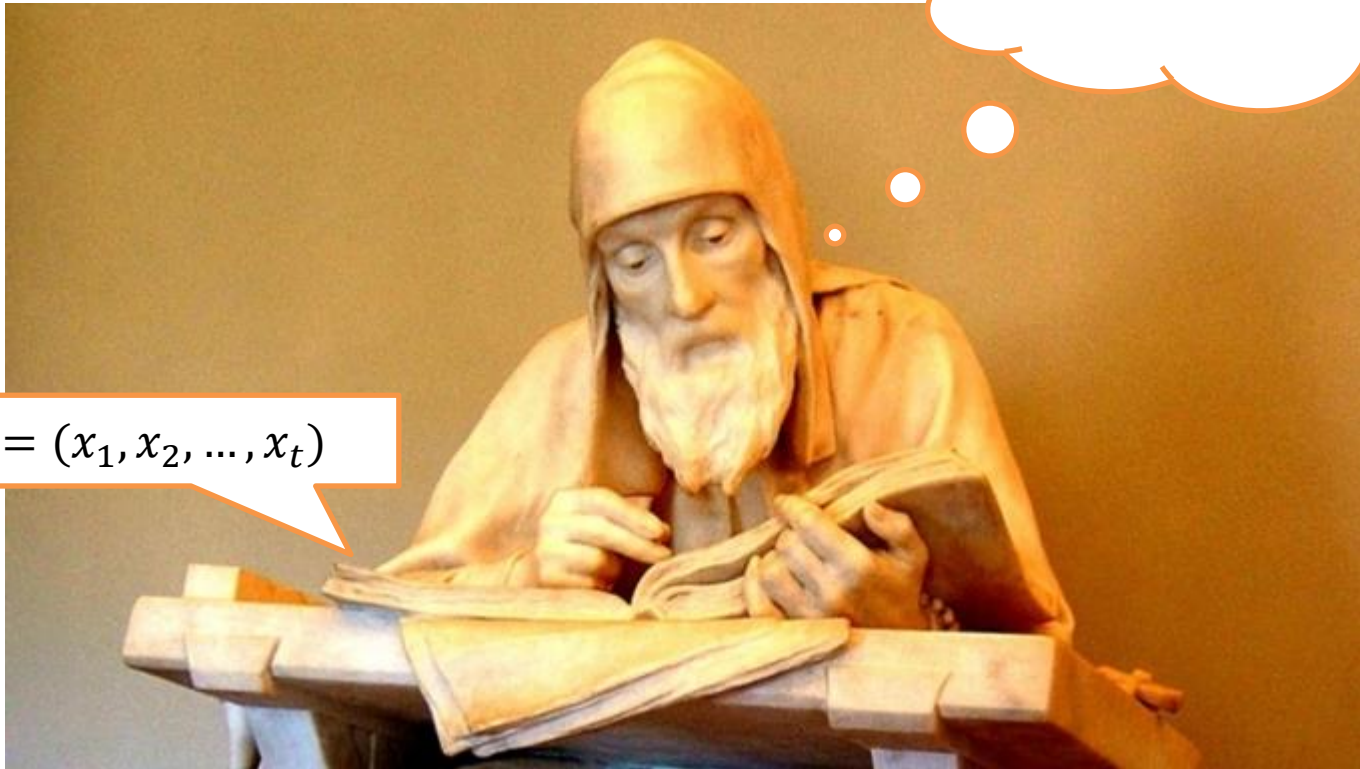
<u>wc_lyricsnet_artists</u>	
id	PK
lyricsnet_id	
lyricsnet_tag	
name	
albums_count	
songs_count	
parsed	

<u>wc_lyricsnet_songs</u>	
id	PK
lyricsnet_id	
artist_id	
album_id	
title	
parsed	
views_count	
view_last_time	
skipped	
archived	

Noisy channel model

$$w = (w_1, w_2, \dots, w_s)$$

$$x = (x_1, x_2, \dots, x_t)$$



$$\hat{w} = \operatorname{argmax}_{w \in V} P(w|x) = \operatorname{argmax}_{w \in V} \frac{P(x|w)P(w)}{P(x)}$$

$$\hat{w} = \operatorname{argmax}_{w \in V} P(x|w)P(w)$$

- $P(x|w)$ – *likelihood(channel model)*
- $P(w)$ – *prior probability*

$$\hat{w} = \operatorname{argmax}_{w \in V} \log P(w) + \log P(x|w)$$

$$P(x|w) = P(f|w \rightarrow x) = P(f) = \prod_{i=1}^k P(f_i)$$

$$\hat{w} = \operatorname{argmax}_{w \in V} \log P(w) + \sum_i \log P(f_i)$$

$$P(x|w) = \begin{cases} \frac{\text{del}[x_{i-1}, w_i]}{\text{count}[x_{i-1} w_i]}, & \text{if deletion} \\ \frac{\text{ins}[x_{i-1}, w_i]}{\text{count}[w_{i-1}]}, & \text{if insertion} \\ \frac{\text{sub}[x_i, w_i]}{\text{count}[w_i]}, & \text{if substitution} \\ \frac{\text{trans}[w_i, w_{i+1}]}{\text{count}[w_i w_{i+1}]}, & \text{if transposition} \end{cases}$$

Модифицированный алгоритм Дамерау-Левенштейна

Базовые операции:

- ADD
- DELETE
- SUBSTITUTE
- TRANSPOSE

Дополнительные операции:

- DOUBLE
- SINGLE

Сбор статистики об ошибках реальных пользователей

Типичные ошибки и опечатки

Данный опрос посвящен выявлению наиболее вероятных ошибок при наборе англоязычных текстов (в частности, поисковых запросов в социальных сетях при поиске музыки и сообщений).

** Обязательно*

Насколько часто при печати вы набираете лишние буквы (символы) в английских словах (например, transposition -> transposditiion) ***

12345

Очень редко☐☐☐☐☐Очень часто

Насколько часто при печати вы пропускаете нужные буквы в английских словах (например, transposition -> transpsition) ***

12345

Очень редко☐☐☐☐☐Очень часто

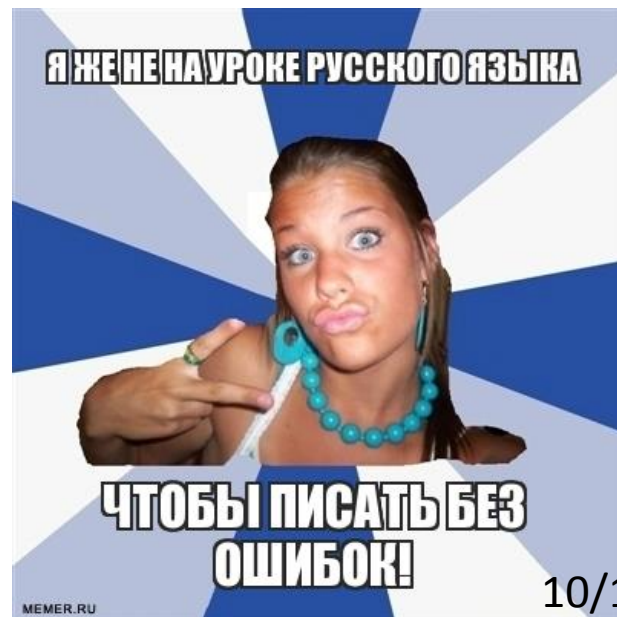
add	delete	transpose	single/double	substitute
2,90	2,90	2,99	2,62	2,83

Немного статистики

Согласно лингвистическим исследованиям:

- 80% всех ошибочных написаний слова находятся на редакционном расстоянии, равном одному от его верного написания;
- Почти в 100% случаев редакционное расстояние равно двум

P.S. Некоторые индивиды, конечно, ошибаются чаще ;)

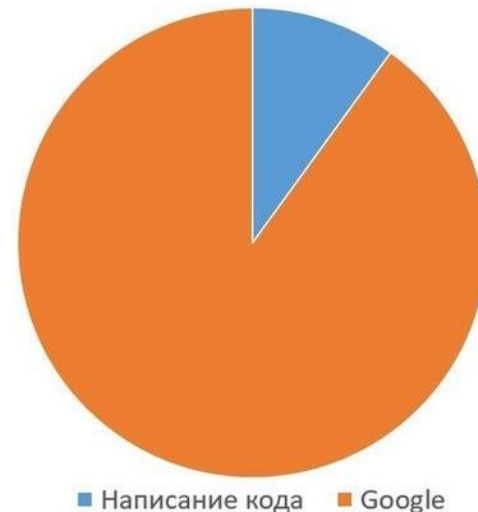


Этапы обработки запросов:

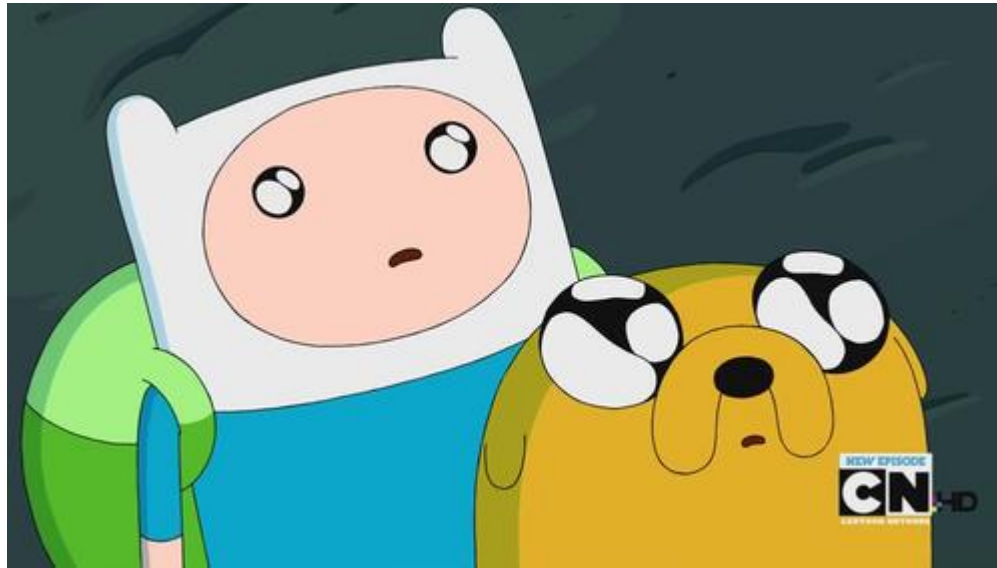
- Разбиваем поисковый запрос на слова
- Для каждого слова применяем алгоритм Noisy Channel и находим список основных песен-кандидатов
- Пересекаем списки для всех слов в запросе, получаем суженный набор кандидатов
- Выбираем из данного набора оптимальных кандидатов, используя метрику Левенштейна



Распределение времени
при программировании



It's demonstration time!

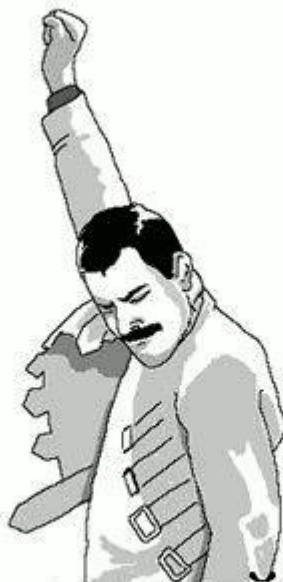


ИСТОЧНИКИ

- Daniel Jurafsky & James H. Martin «Speech and language processing»
- Kernighan, M. D., Church, K.W., and Gale, W. A. (1990). A spelling correction program based on a noisy channel model.
- Norvig, P. (2009). Natural language corpus data (<http://norvig.com/ngrams/>)

Спасибо за внимание!

Просто так взяло и



СКОМПИЛИЛОСЬ